# A Model-based Trust Region Method for Stochastic Derivative-free Optimization

Jeffrey Larson    Stephen Billups

**Argonne National Laboratory**

July 26, 2015

# The Problem

We want to solve:

$$\underset{x \in \mathbb{R}^n}{\text{minimize}}\, f(x)$$

when $\nabla f(x)$ is unavailable and we only have access to noise-corrupted function evaluations $\bar{f}(x)$.

## The Problem

We want to solve:

$$\underset{x \in \mathbb{R}^n}{\text{minimize}} \, f(x)$$

when $\nabla f(x)$ is unavailable and we only have access to noise-corrupted function evaluations $\bar{f}(x)$.

Such noise may be deterministic (e.g., from iterative methods) or stochastic (e.g., from a Monte-Carlo process).

## The Problem

We want to solve:

$$\underset{x \in \mathbb{R}^n}{\text{minimize}}\, f(x)$$

when $\nabla f(x)$ is unavailable and we only have access to noise-corrupted function evaluations $\bar{f}(x)$.

Such noise may be deterministic (e.g., from iterative methods) or stochastic (e.g., from a Monte-Carlo process).

Model-based methods are one of the most popular methods when $\nabla f$ is unavailable, and the only recourse when noise is deterministic.

## The Problem

We want to solve:

$$\underset{x \in \mathbb{R}^n}{\text{minimize}}\, f(x)$$

when $\nabla f(x)$ is unavailable and we only have access to noise-corrupted function evaluations $\bar{f}(x)$.

Such noise may be deterministic (e.g., from iterative methods) or stochastic (e.g., from a Monte-Carlo process).

Model-based methods are one of the most popular methods when $\nabla f$ is unavailable, and the only recourse when noise is deterministic.

$n$ is small, $f$ is likely nonconvex.

## The Problem

We analyze the convergence of our method in the stochastic case:

$$\bar{f}(x) = f(x) + \epsilon,$$

where $\epsilon$ is identically distributed with mean 0 and variance $\sigma^2 < \infty$.

## The Problem

We analyze the convergence of our method in the stochastic case:
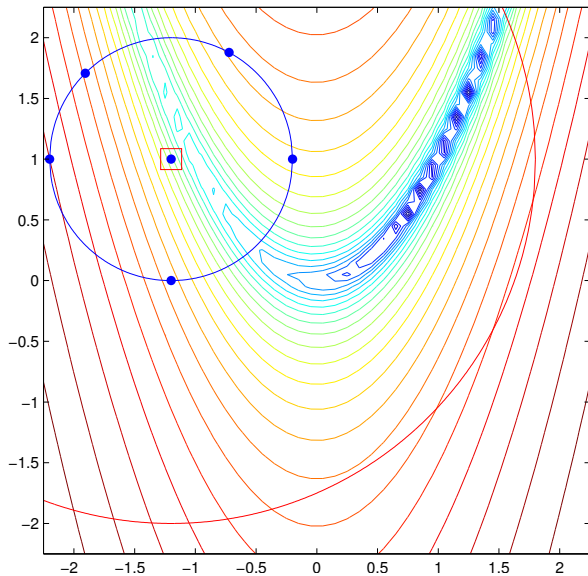
$$\bar{f}(x) = f(x) + \epsilon,$$

where $\epsilon$ is identically distributed with mean 0 and variance $\sigma^2 < \infty$.
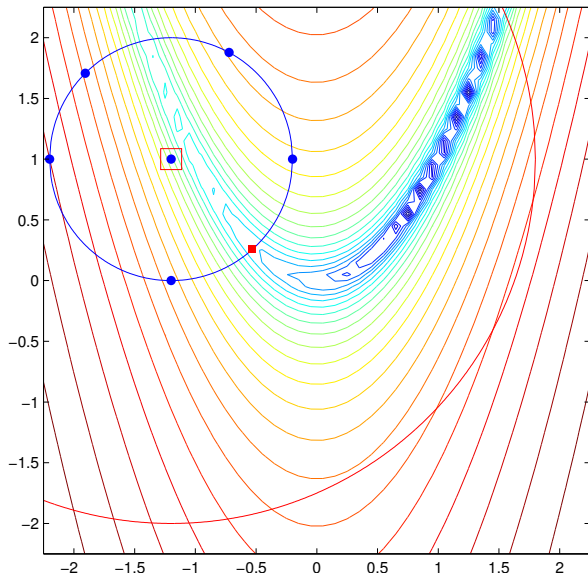
This is equivalent to solving:

$$\underset{x}{\text{minimize}} \ \mathbb{E}\left[\bar{f}(x)\right].$$

# Prototype

# Prototype

# Prototype

# Prototype

# Prototype

# Prototype

# Prototype

# Prototype

# Prototype

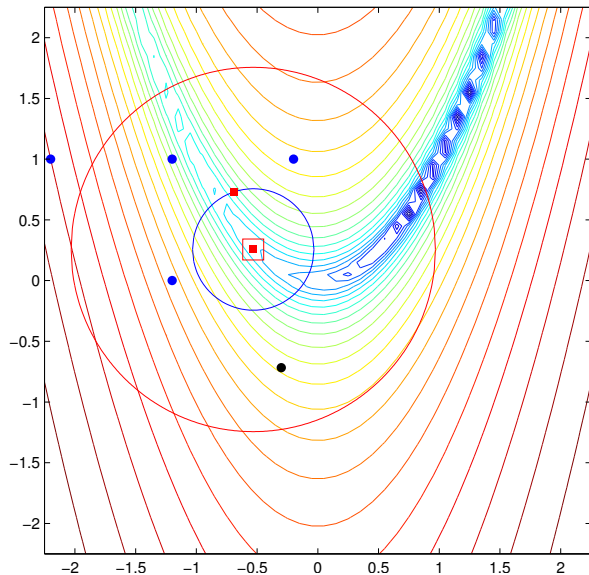# Prototype

# Prototype

# Prototype

# Prototype

# Prototype

# Prototype

# Strongly Λ-poised Sets

# Stochastic Approximation

Iterates usually have the form:

$$x^{k+1} = x^k + a_k \, G(x^k),$$

where

- $G(x^k)$ is a cheap, unbiased estimate for $\nabla f(x^k)$

# Stochastic Approximation

Iterates usually have the form:

$$x^{k+1} = x^k + a_k G(x^k),$$

where

- $G(x^k)$ is a cheap, unbiased estimate for $\nabla f(x^k)$

  - For Kiefer-Wolfowitz,

    $$G_i(x^k) = \frac{\bar{f}(x^k + c_k e_i) - \bar{f}(x^k - c_k e_i)}{2 c_k}$$

    where $e_i$ is the $i$th column of $I_n$.

# Stochastic Approximation

Iterates usually have the form:

$$x^{k+1} = x^k + a_k G(x^k),$$

where

- $G(x^k)$ is a cheap, unbiased estimate for $\nabla f(x^k)$

  - For Spall's SPSA,

    $$G_i(x^k) = \frac{\bar{f}(x^k + c_k \delta^k) - \bar{f}(x^k - c_k \delta^k)}{2 c_k \delta_i^k}$$

    where $\delta^k \in \mathbb{R}^n$ is a random perturbation vector

# Stochastic Approximation

Iterates usually have the form:

$$x^{k+1} = x^k + a_k\, G(x^k),$$

where

- $G(x^k)$ is a cheap, unbiased estimate for $\nabla f(x^k)$

- $a_k$ is a sequence of step sizes

# Stochastic Approximation

Iterates usually have the form:

$$x^{k+1} = x^k + a_k G(x^k),$$

where

- $G(x^k)$ is a cheap, unbiased estimate for $\nabla f(x^k)$

- $a_k$ is a sequence of step sizes (specified by the user) satisfying:

$$\sum_{k=1}^{\infty} a_k = \infty \qquad \lim_{k \to \infty} a_k = 0$$

## Stochastic Approximation

Iterates usually have the form:

$$x^{k+1} = x^k + a_k G(x^k),$$

where

- $G(x^k)$ is a cheap, unbiased estimate for $\nabla f(x^k)$

- $a_k$ is a sequence of step sizes (specified by the user) satisfying:

$$\sum_{k=1}^{\infty} a_k = \infty \qquad \lim_{k \to \infty} a_k = 0$$

Algorithm performance depends significantly on sequence $a_k$.

# Response Surface Methodology

- Build models using a fixed pattern of points, (e.g., cubic, spherical, or orthogonal designs).

- Finding the design that constructs response surfaces approximating the function (without few function evaluations) can be difficult for problems where the user has no prior expertise.

# Modifications to Existing Methods

Take a favorite method and repeatedly evaluate the function at points of interest.

- ▶ Stochastic approximation modified by Dupuis, Simha (1991)
- ▶ Response surface methods modified by Chang et al. (2012)
- ▶ UOBYQA modified by Deng, Ferris (2006)
- ▶ Nelder-Mead modified by Tomick et al. (1995)
- ▶ DIRECT modified by Deng, Ferris (2007)

# Modifications to Existing Methods

Take a favorite method and repeatedly evaluate the function at points of interest.

- ▶ Stochastic approximation modified by Dupuis, Simha (1991)
- ▶ Response surface methods modified by Chang et al. (2012)
- ▶ UOBYQA modified by Deng, Ferris (2006)
- ▶ Nelder-Mead modified by Tomick et al. (1995)
- ▶ DIRECT modified by Deng, Ferris (2007)

There are two downsides to such an approach:

1. Repeated sampling provides information about the noise $\epsilon$, not $f$.
2. If the noise is deterministic, no information is gained.

# Overview

We therefore desire a method that

1. Adjusts the step size as it progresses

2. Does not use a fixed design of points

3. Does not repeatedly sample points

# Overview

We therefore desire a method that

1. Adjusts the step size as it progresses

2. Does not use a fixed design of points

3. Does not repeatedly sample points

We'd like the class of possible models to be general.

# $\kappa$-fully Linear model

## Definition

If $f \in LC$ and $\exists$ a vector $\kappa = (\kappa_{ef}, \kappa_{eg})$ of positive constants such that

- the error between the gradient of the model and the gradient of the function satisfies

$$\|\nabla f(y) - \nabla m(y)\| \leq \kappa_{eg}\Delta \ \ \forall y \in B(x; \Delta),$$

- the error between the model and the function satisfies

$$|f(y) - m(y)| \leq \kappa_{ef}\Delta^2 \ \ \forall y \in B(x; \Delta),$$

we say the model is $\kappa$-*fully linear* on $B(x; \Delta)$.

# $\alpha$-probabilistically $\kappa$-fully Linear model

## Definition

Let $\kappa = (\kappa_{ef}, \kappa_{eg})$ be a given vector of constants, and let $\alpha \in (0, 1)$. Let $B \subset \mathbb{R}^n$ be given. A random model $m_k$ generated at the $k$th iteration of an algorithm is $\alpha$-probabilistically $\kappa$-fully linear on $B$ if

$$P\left(m_k \text{ is a } \kappa\text{-fully linear model of } f \text{ on } B \middle| \mathcal{F}_{k-1}\right) \geq \alpha,$$

where $\mathcal{F}_{k-1}$ denotes the realizations of all the random events for the first $k - 1$ iterations.

# Regression Models can be $\alpha$-probabilistically $\kappa$-fully Linear

## Theorem

*For a given $x \in \mathbb{R}^n$, $\Delta > 0$, $\alpha \in (0, 1)$,*

- *$Y \subset B(x; \Delta)$ is strongly $\Lambda$-poised,*
- *The noise present in $\bar{f}$ is i.i.d. with mean 0, variance $\sigma^2 < \infty$,*
- *$|Y| \geq C/\Delta^4$,*

*Then there exist constants $\kappa = (\kappa_{ef}, \kappa_{eg})$ (independent of $\Delta$ and $Y$) such that the linear model $m$ regressing $Y$ is $\alpha$-probabilistically $\kappa$-fully linear on $B(x; \Delta)$.*

# Measuring Progress

In traditional trust region methods, if $x^k + s^k$ is the minimizer of $m_k$, the success of moving from $x^k$ to $x^k + s^k$ is measured by

$$\rho_k = \frac{f(x^k) - f(x^k + s^k)}{m_k(x^k) - m_k(x^k + s^k)}$$

# Measuring Progress

In traditional trust region methods, if $x^k + s^k$ is the minimizer of $m_k$, the success of moving from $x^k$ to $x^k + s^k$ is measured by

$$\rho_k = \frac{f(x^k) - f(x^k + s^k)}{m_k(x^k) - m_k(x^k + s^k)}$$

In the stochastic case, a similar calculation is not obvious.

$$\rho_k = \frac{\bar{f}(x^k) - \bar{f}(x^k + s^k)}{m_k(x^k) - m_k(x^k + s^k)}$$

# Measuring Progress

In traditional trust region methods, if $x^k + s^k$ is the minimizer of $m_k$, the success of moving from $x^k$ to $x^k + s^k$ is measured by

$$\rho_k = \frac{f(x^k) - f(x^k + s^k)}{m_k(x^k) - m_k(x^k + s^k)}$$

In the stochastic case, a similar calculation is not obvious.

$$\rho_k = \frac{m_k(x^k) - m_k(x^k + s^k)}{m_k(x^k) - m_k(x^k + s^k)}$$

# Measuring Progress

In traditional trust region methods, if $x^k + s^k$ is the minimizer of $m_k$, the success of moving from $x^k$ to $x^k + s^k$ is measured by

$$\rho_k = \frac{f(x^k) - f(x^k + s^k)}{m_k(x^k) - m_k(x^k + s^k)}$$

In the stochastic case, a similar calculation is not obvious.

$$\rho_k = \frac{m_k(x^k) - \hat{m}_k(x^k + s^k)}{m_k(x^k) - m_k(x^k + s^k)}$$

# Measuring Progress

In traditional trust region methods, if $x^k + s^k$ is the minimizer of $m_k$, the success of moving from $x^k$ to $x^k + s^k$ is measured by

$$\rho_k = \frac{f(x^k) - f(x^k + s^k)}{m_k(x^k) - m_k(x^k + s^k)}$$

In the stochastic case, a similar calculation is not obvious.

$$\rho_k = \frac{F_k^0 - F_k^s}{m_k(x^k) - m_k(x^k + s^k)}$$

# One Last Part

For our analysis, we need estimates of $f(x^k)$ and $f(x^k + s^k)$ that are slightly different than those provided by the model functions.

Let $F_k^0$ and $F_k^s$ denote the sequence of estimates of $f(x^k)$ and $f(x^k + s^k)$.

We need to be able to construct estimates satisfying

$$\mathbb{P}\left[\left|F_k^0 - f(x^k)\right| > \epsilon \min\left\{\Delta_k, \Delta_k^2\right\} \big| \mathcal{F}_{k-1}\right] < \theta$$

$$\text{and } \mathbb{P}\left[\left|F_k^s - f(x^k + s^k)\right| > \epsilon \min\left\{\Delta_k, \Delta_k^2\right\} \Big| \mathcal{F}_{k-1}\right] < \theta,$$

for any $\epsilon > 0$ and $\theta > 0$.

**Algorithm 1:** A trust region algorithm to minimize a stochastic function

---

Set $k = 0$;

**Start**

Build a $\alpha$-probabilistically $\kappa$-fully linear model $m_k$ on $B(x^k; \Delta_k)$;

Compute $s^k = \arg \min\limits_{s: \|x^k - s\| \leq \Delta_k} m_k(s)$;

**if** $m_k(s^k) - m_k(x^k + s^k) \geq \beta \Delta_k$ **then**

    Calculate $\rho_k = \dfrac{F_k^0 - F_k^s}{m_k(x^k) - m_k(x^k + s^k)}$;

    **if** $\rho_k \geq \eta$ **then**

        Calculate $x^{k+1} = x^k + s^k$; $\Delta_{k+1} = \gamma_{inc} \Delta_k$;

    **else**

        $x^{k+1} = x^k$; $\Delta_{k+1} = \gamma_{dec} \Delta_k$;

    **end**

**else**

    $x^{k+1} = x^k$; $\Delta_{k+1} = \gamma_{dec} \Delta_k$;

**end**

$k = k + 1$ and go to **Start**;

---

# Convergence

Under what assumptions will our algorithm converge almost surely to a first-order stationary point?

- Assumptions on $f$

- Assumptions on $\epsilon$

- Assumptions on algorithmic constants

# Convergence

### Assumption

*On some set $\Omega \subseteq \mathbb{R}^n$ containing all iterates visited by the algorithm,*

- $\nabla f$ *is Lipschitz continuous with constant $L_g$*
- *$f$ has bounded level sets*

### Assumption

*The additive noise $\epsilon$ observed when computing $\bar{f}$ is independent and identically distributed with mean zero and bounded variance $\sigma^2$.*

# Convergence

### Assumption

The constants $\alpha \in (0, 1)$, $\gamma_{dec} \in (0, 1)$, and $\gamma_{inc} > 1$ **satisfy**

$$\alpha \geq \max \left\{ \frac{1}{2}, 1 - \frac{\frac{\gamma_{inc} - 1}{\gamma_{inc}}}{4 \left[ \frac{\gamma_{inc} - 1}{2\gamma_{inc}} + \frac{1 - \gamma_{dec}}{\gamma_{dec}} \right]} \right\},$$

where

- $\alpha$ is the lower bound on the probability of having a $\kappa$-fully linear model,
- $\gamma_{dec} \in (0, 1)$ is the factor by which we decrease the trust region radius,
- $\gamma_{inc} > 1$ is the factor by which the trust radius is increased.

# Convergence

## Assumption

*The constants $\alpha \in (0, 1)$, $\gamma_{dec} \in (0, 1)$, and $\gamma_{inc} > 1$ satisfy*

$$\alpha \geq \max \left\{ \frac{1}{2}, 1 - \frac{\frac{\gamma_{inc}-1}{\gamma_{inc}}}{4 \left[ \frac{\gamma_{inc}-1}{2\gamma_{inc}} + \frac{1-\gamma_{dec}}{\gamma_{dec}} \right]} \right\},$$

*where*

- $\alpha$ *is the lower bound on the probability of having a $\kappa$-fully linear model,*
- $\gamma_{dec} \in (0, 1)$ *is the factor by which we decrease the trust region radius,*
- $\gamma_{inc} > 1$ *is the factor by which the trust radius is increased.*

If $\gamma_{inc} = 2$ and $\gamma_{dec} = 0.5 \implies \alpha \geq 0.9$.
If $\gamma_{inc} = 2$ and $\gamma_{dec} = 0.9 \implies \alpha \geq 0.65$.

# Proof Outline

## Theorem

*If the above assumptions are satisfied, our algorithm converges almost surely to a first-order stationary point of $f$.*

- Show the sequence of trust region radii $\Delta_k \to 0$ almost surely.

# Proof Outline

## Theorem

*If the above assumptions are satisfied, our algorithm converges almost surely to a first-order stationary point of $f$.*

- Show the sequence of trust region radii $\Delta_k \to 0$ almost surely.
- Show if $\Delta_k$ ever falls below some constant multiple of the model gradient, $\Delta_{k+1} > \Delta_k$ with high probability.

# Proof Outline

## Theorem

*If the above assumptions are satisfied, our algorithm converges almost surely to a first-order stationary point of $f$.*

- Show the sequence of trust region radii $\Delta_k \to 0$ almost surely.
- Show if $\Delta_k$ ever falls below some constant multiple of the model gradient, $\Delta_{k+1} > \Delta_k$ with high probability.
- Lastly, show that, the sequence of ratios

$$\{\psi_k\} = \left\{ \frac{\left\| \nabla f(x^k) \right\|}{\Delta_k} \right\}$$

satisfies $\mathbb{E}\left[\psi_{k+1} | \mathcal{F}_k\right] \leq \psi_k$ when $\psi_k \geq L$. This allows us to prove $\left\| \nabla f(x^k) \right\| \to 0$ in probability.

**Algorithm 1:** A trust region algorithm to minimize a stochastic function

---

Set $k = 0$;

**Start**

Build a $\alpha$-probabilistically $\kappa$-fully linear model $m_k$ on $B(x^k; \Delta_k)$;

Compute $s^k = \arg \min\limits_{s: \|x^k - s\| \leq \Delta_k} m_k(s)$;

**if** $m_k(s^k) - m_k(x^k + s^k) \geq \beta \Delta_k$ **then**

    Calculate $\rho_k = \dfrac{F_k^0 - F_k^s}{m_k(x^k) - m_k(x^k + s^k)}$;

    **if** $\rho_k \geq \eta$ **then**

        Calculate $x^{k+1} = x^k + s^k$; $\Delta_{k+1} = \gamma_{inc}\Delta_k$;

    **else**

        $x^{k+1} = x^k$; $\Delta_{k+1} = \gamma_{dec}\Delta_k$;

    **end**

**else**

    $x^{k+1} = x^k$; $\Delta_{k+1} = \gamma_{dec}\Delta_k$;

**end**

$k = k + 1$ and go to **Start**;

---

## Prototype

- $m_k$ is a linear regression model on a sample set of $(n+1)C_k$ sample points, where $C_k$ is defined by

$$C_k = \left\lceil \frac{k}{1000} \right\rceil \frac{\max\left\{n+1, \left\lfloor \frac{1}{\Delta_k^4} \right\rfloor\right\}}{n+1}.$$

The sample set consists of $C_k$ randomly rotated copies of the set

$$\{x^k, x^k + \Delta_k e_1, \ldots, x^k + \Delta_k e_n\}$$

## Prototype

- $m_k$ is a linear regression model on a sample set of $(n+1)C_k$ sample points, where $C_k$ is defined by

$$C_k = \left\lceil \frac{k}{1000} \right\rceil \frac{\max\left\{ n+1, \left\lfloor \frac{1}{\Delta_k^4} \right\rfloor \right\}}{n+1}.$$

The sample set consists of $C_k$ randomly rotated copies of the set

$$\{x^k, x^k + \Delta_k e_1, \ldots, x^k + \Delta_k e_n\}$$

- $F_k^0 = m_k^0(x^k)$, where $m_k^0$ is a linear regression model using $C_k$ randomly rotated copies of the set

$$\{x^k, x^k + 0.5\Delta_k e_1, \ldots, x^k + 0.5\Delta_k e_n\}$$

# Prototype

- $m_k$ is a linear regression model on a sample set of $(n+1)C_k$ sample points, where $C_k$ is defined by

$$C_k = \left\lceil \frac{k}{1000} \right\rceil \frac{\max\left\{ n+1, \left\lfloor \frac{1}{\Delta_k^4} \right\rfloor \right\}}{n+1}.$$

The sample set consists of $C_k$ randomly rotated copies of the set

$$\{x^k, x^k + \Delta_k e_1, \ldots, x^k + \Delta_k e_n\}$$

- $F_k^s = m_k^s(x^k)$, where $m_k^s$ is a linear regression model using $C_k$ randomly rotated copies of the set

$$\{x^k + s^k, x^k + s^k + 0.5\Delta_k e_1, \ldots, x^k + s^k + 0.5\Delta_k e_n\}$$

# Problem Set

53 problems of the form:

$$f(x) = \sum_{i=1}^{m} \left[ (1 + \sigma) F_i(x) \right]^2,$$

where $\sigma \sim U[-0.1, 0.1]$.

## Problem Set

53 problems of the form:

$$f(x) = \sum_{i=1}^{m} [(1 + \sigma)F_i(x)]^2 \,,$$

where $\sigma \sim U[-0.1, 0.1]$.

If $S$ is the set of solvers to be compared on a suite of problems $P$, let $t_{p,s}$ be the number of iterates required for solver $s \in S$ on a problem $p \in P$ to find a function value satisfying:

$$f(x) - f_L \leq \tau \left( f(x^0) - f_L \right),$$

where $f_L$ is the best function value achieved by any $s \in S$.

# Problem Set

### Comments

- We are using the true function value $f$, not the observed $\bar{f}$.
- Since the noise is stochastic, each solver is run 10 times per problem.

If $S$ is the set of solvers to be compared on a suite of problems $P$, let $t_{p,s}$ be the number of iterates required for solver $s \in S$ on a problem $p \in P$ to find a function value satisfying:

$$f(x) - f_L \leq \tau \left( f(x^0) - f_L \right),$$

where $f_L$ is the best function value achieved by any $s \in S$.

# Performance Profile

Then the performance profile of a solver $s \in S$ is the following fraction:

$$\rho_s(\phi) = \frac{1}{|P|} \left| \left\{ p \in P : \frac{t_{p,s}}{\min\left\{ t_{p,s} : s \in S \right\}} \leq \phi \right\} \right|$$

# Performance Profile

Then the performance profile of a solver $s \in S$ is the following fraction:

$$\rho_s(\phi) = \frac{1}{|P|} \left| \left\{ p \in P : \frac{t_{p,s}}{\min\{t_{p,s} : s \in S\}} \leq \phi \right\} \right|$$

- $\rho_s(1)$: Fraction of $P$ method $s$ solves first.
- $\lim_{\phi \to \infty} \rho_s(\phi)$: Fraction of $P$ method $s$ eventually solves.
- $\rho_s(\phi)$: Fraction of $P$ method $s$ solves in under $\phi$ times the evaluations required for the best method.
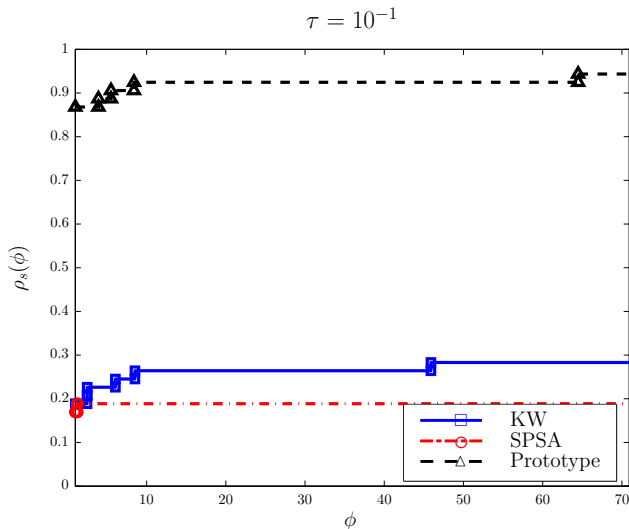
# Performance Profile

We compare our prototype against Spall's versions of Kiefer-Wolfowitz and SPSA with step sizes as recommended in Sections 6.6 and 7.5.2 of Spall (2003)

$$a_k = \frac{1}{(k+1+A)^{0.602}} \qquad c_k = \frac{1}{(k+1)^{0.101}}$$

where $A$ is one tenth of the total budget of function evaluations.

# Performance Profile
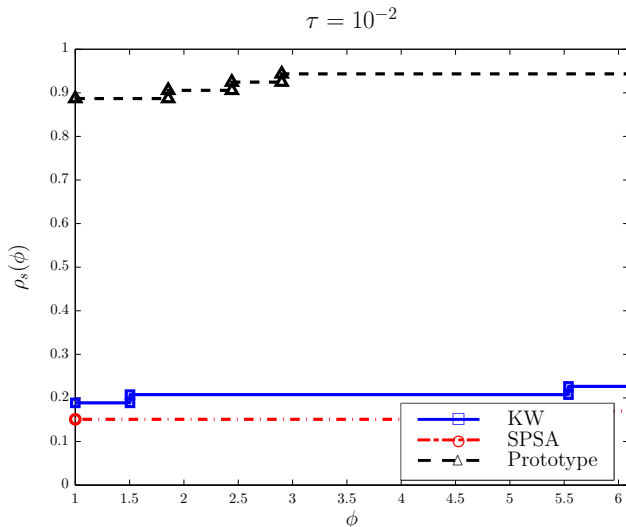


$\tau = 10^{-1}$

$\rho_s(1)$:
Fraction s solves first

$\lim_{\phi \to \infty} \rho_s(\phi)$:
Fraction s solves

$\rho_s(\phi)$:
Fraction s solves in under $\phi$ times the evaluations required for the best method.

Legend: KW, SPSA, Prototype

# Performance Profile



$\rho_s(1)$:
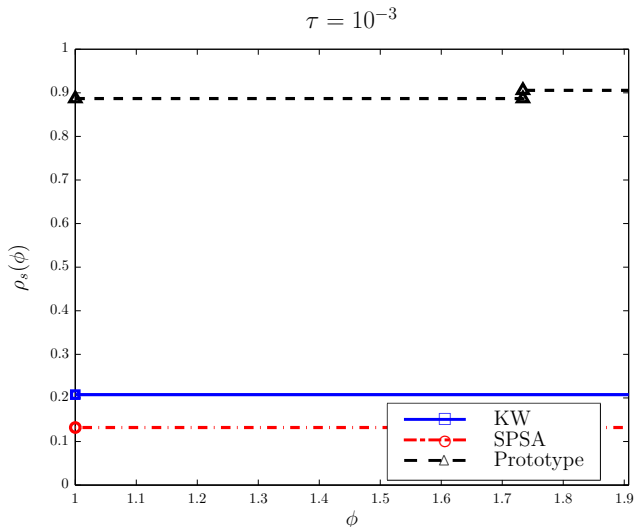Fraction $s$ solves first

$\lim_{\phi \to \infty} \rho_s(\phi)$:
Fraction $s$ solves

$\rho_s(\phi)$:
Fraction $s$ solves in under $\phi$ times the evaluations required for the best method.

# Performance Profile



$\tau = 10^{-3}$

$\rho_s(1)$:
Fraction $s$ solves first

$\lim_{\phi \to \infty} \rho_s(\phi)$:
Fraction $s$ solves

$\rho_s(\phi)$:
Fraction $s$ solves in under $\phi$ times the evaluations required for the best method.
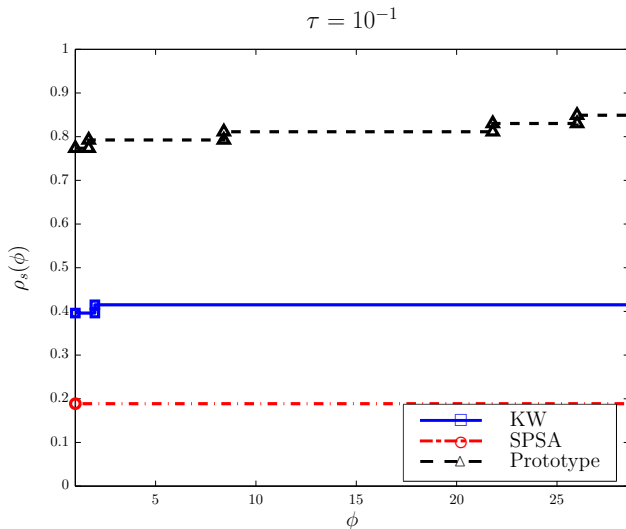
## Another Problem Set

53 problems of the form:

$$f(x) = \sigma_p + \sum_{i=1}^{m} [F_i(x)]^2,$$

where $\sigma_p \sim N\left(0, (0.1\Delta_p)^2\right)$ and $\Delta_p = \sum_i F_i(x^0) - \sum_i F_i(x^*)$.

# Performance Profile



$\tau = 10^{-1}$

$\rho_s(1)$:
Fraction $s$ solves first
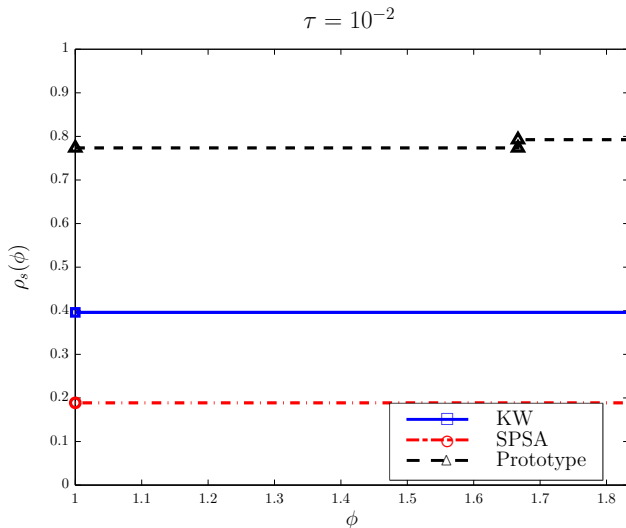
$\lim_{\phi \to \infty} \rho_s(\phi)$:
Fraction $s$ solves

$\rho_s(\phi)$:
Fraction $s$ solves in under $\phi$ times the evaluations required for the best method.

Legend:
- KW
- SPSA
- Prototype

## Performance Profile



$\tau = 10^{-2}$

$\rho_s(1)$:
Fraction $s$ solves first

$\lim_{\phi \to \infty} \rho_s(\phi)$:
Fraction $s$ solves

$\rho_s(\phi)$:
Fraction $s$ solves in under $\phi$ times the evaluations required for the best method.

# Further Information and Current Work

"Stochastic Derivative-free Optimization using a Trust Region Framework"

# Further Information and Current Work

## Preprint on Optimization Online

"Stochastic Derivative-free Optimization using a Trust Region Framework"

- ▶ Generalizing results to ensure a practical algorithm converges.
  - ▶ For example, not requiring $\alpha$-probabilistically $\kappa$-fully linear models every iteration.

# Further Information and Current Work

"Stochastic Derivative-free Optimization using a Trust Region Framework"

- ▶ Generalizing results to ensure a practical algorithm converges.
    - ▶ For example, not requiring $\alpha$-probabilistically $\kappa$-fully linear models every iteration.

- ▶ Smartly constructing $\alpha$-probabilistically $\kappa$-fully linear models.